



THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer

Reducing OCR errors by combining two OCR systems

Citation for published version:

Volk, M, Marek, T & Sennrich, R 2010, Reducing OCR errors by combining two OCR systems. in C Sporleder & K Zervanou (eds), *Proceedings of the ECAI 2010 Workshop on Language Technology for Cultural Heritage, Social Sciences, and Humanities (LaTeCH 2010)*. pp. 61-65.
<<http://dx.doi.org/10.5167/uzh-35259>>

Link:

[Link to publication record in Edinburgh Research Explorer](#)

Document Version:

Peer reviewed version

Published In:

Proceedings of the ECAI 2010 Workshop on Language Technology for Cultural Heritage, Social Sciences, and Humanities (LaTeCH 2010)

General rights

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact openaccess@ed.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.



Reducing OCR Errors by Combining Two OCR Systems.

Martin Volk, Torsten Marek and Rico Sennrich

Institute of Computational Linguistics
University of Zurich

16 August 2010

- Corpus of "Alpine" texts
- Mountaineering and travel reports in multiple languages
- For studies on language development, culture, technology, ...
- For geo-tagging, machine translation, ...
- First step: Publications of the Swiss Alpine Club since 1864



Aus der Landschaft Davos.

251

streifen, der das ganze Thal überzieht, steht in gleicher Meereshöhe eine stündig bewohnte Hütte, worin eine alte, tapfere Jungfrau mit ihren Ziegen, Katzen und Hühnern lauschalet; und willig trotzt sie allen Winterstürmen und lauscht, ob der Frühling wiederkehre, der sie hervorlocke zur Bestellung ihres Ackerleins. Keine besondere Verlassenheit, kein Mangel zwingt die resolute Person zu ihrer Einsamkeit, und wer sie fragt, ob sie sich darin nicht unglücklich fühle, wird nicht Übel heimgeschickt. So viel und so wenig braucht der Mensch manchemal, um zufrieden zu sein!

Nur wenige menschliche Wohnungen steigen in Davos noch höher ins Gebirge, und locken zum Ausharren auch im Winter in einer Region, die keinen Ackerbau mehr kennt, und nur der Viehzucht und einem Gemüsegärtchen noch eine Stütze gewährt. Auf der Schatzalp ob Davosplatz,



Bärenhaus in der Kuppe.

Photographie von J. Bille in Davos.

an der obern Waldgrenze (1875^o), hat der Fremdenbesuch auch im Winter zu ständiger Besiedelung geführt; im Filletal ist es die Paiststraße, die in dem wilden Hochgebirgsthale noch auf der „Alpenrose“ bei 1830 Metern und auf dem Tschuggen bei 1941 Metern den Menschen dem rauhen Klima trotzen liess. Das Filletalshospiz, den grössten Teil des Jahres hindurch in arktischer Umgebung (2388^o), liegt schon auf Engadinerboden. Das Diechmalthal, an den Gletschern der Scaloetta endigend, ist bald des Waldes entbehrend, ist zur Winterzeit nur im vordern Thal bewohnt, und die letzten, ständig besetzten Häuser stehen schon bei 1700 Metern. Viel späterlicher bewohnt als das fremde Sertig, kennen Flüela und Dischma kaum mehr den Ackerbau.

Schon zeitig im Frühling, viel eher, als wohl der Tiefkinder glaubt, der bis in den Sommer die Alpengegenden noch unter hohen Schne-

**Sicherheit,
Medizin,
Rettungswesen**
**Sicurezza, medicina,
soccorso in montagna**
**Sécurité, médecine,
sauvetage**

Zurück in der Natur, als
Hütte im Wald und nicht
auf der Strasse. Diese Hütte
ist ein altes Haus, das aus
einem Steinhaus (Hütte)
aus dem 18. Jahrhundert
stammt.

Wenn man diese Hütte
mit der Hütte auf der
Strasse vergleicht, so
kann man sehen, dass
die Hütte auf der Strasse
nicht mehr existiert.

Zecken in den Klettergärten

Die Blutsauger klettern immer höher

Ein Teil der Zecken in der Schweiz
übertragen Krankheiten, unter an-
deren das FSME-Virus, das eine
schwere Hirnhautentzündung auslö-
sen kann. Eine Studie belegt nun,
dass FSME-verseuchte Zecken auch
in Höhenlagen vorkommen, in denen
sie bisher nicht vermutet worden
sind.

Die milden Temperaturen locken Alpi-
nisten, Sportkletterern und Wanderer
wieder vor die Tür. Doch die Wärme, die
zum Tragen kurzer und leichter Klei-
dung verleitet, hat auch ihre Tücken. So
lausern etwa Zecken in freudigem Unter-
holz, an Waldstamm und in Hecken auf
ihre Opfer. Hier haben die kleinen, teil-
weise lauffähigen Zecken genügend
Chancen, auf ungegeschützte Haut zu ge-
langen. Die von ihnen bevorzugten Kör-
perstellen sind Kniekehlen, Schenkel-
gelenk, Bauchnabel, Achselhöhlen, Schul-
tern, Nacken und hinter das Ohr.

Nicht nur Wanderer müssen auf-
merksam sein. Die Gefahr, dass die Blutsau-
ger an der Kleidung haften bleiben und
sich auf die Suche nach freiem Kör-
perstellen machen, ist auch im Zeltlager

zu Klettergärten gross, führen die
Wägen doch oft durch niedrigen Gebüsch
und Gras. Hinzu kommt, dass die Tiere
von der Klimaveränderung profitieren
und nun neugieriger in Höhenlagen
von bis zu 1500 Metern zu finden sind.
Bislang fand man Zecken nur bis in Hö-
hen von 1000 Metern. Auch wer sich in

diesen Gebieten aufhält, läuft Gefahr,
mit den Tieren in Kontakt zu kommen,
mildern die Vorlesungen des Labors Spe-
zialisierte.

Suche in aufwendiger Kleinarbeit

Problematisch war bis anhin, dass man
nicht wusste, wo sich die FSME-
befallenen Zecken aufhalten. FSME-
Erkrankungen mussten zwar bisher aus
Bundesamt für Gesundheit (BAG) ge-
meldet werden. Doch leider können die
wenigsten sagen, woher genau sie die
Zecken haben. Zecken mit der Krankheit
erst einige Tage nach dem Zwischenfall
auf. Daher konnten Ärzte oft nur vage
Informationen über den Ort der Infek-
tion aus dem BAG weitergeben.

Eine Studie des Labors Spe-
zialisierte mit Spezialisten der ABC-
Abwehr-Truppe der Armee liefert nun
genauere Daten. Vom April bis zum Juli
2009 wurden an 163 Orten in der Schweiz
62 343 Zecken eingemeldet. «Unter
suchten wurden Gebiete, aus denen Krank-



Zeckenkrankheiten

Zecken in der Schweiz können zwei
Krankheiten mit schweren Folgen über-
tragen: FSME (Frühsommer-Meningo-
enzephalitis) und Lyme-Borreliose.
Das FSME-Virus tritt meist im Frühsommer
auf und wird von Zecken in Heide-
sträucher- und Heidegebieten übertragen.
Eine Impfung ist möglich. Nach einer
Impfung können jedoch nur noch die
Symptome, nicht aber die Krankheit
behandelt werden. Immunität ist
nur bei durchgemachter Krankheit lei-
stungsfähig gewährleistet, die Impfung

hingegen muss alle drei Jahre aufgefrischt
werden.
Die sogenannte Lyme-Borreliose wird von
einem Drittel der Zecken übertragen.
Der Erreger greift Haut, Gelenke, Herzmuskel,
Muskeln und Herz an. Eine Impfung
dagegen ist nicht möglich, jedoch eine Be-
handlung mit Antibiotika nach einem Biss,
wobei danach keine Immunität gegen den
Erreger besteht. Weitere Infos gibt es un-
ter [http://www.labor-spez.ch/de/aktuel-
des.htm](http://www.labor-spez.ch/de/aktuel-
des.htm), unter www.bag.admin.ch, Stich-
wort FSME, oder unter www.zecken.de

- We use commercial OCR systems (Abbyy FineReader and OmniPage)
- OCR systems are language-dependent
- Challenges during OCR of the Text+Berg corpus
 - Corpus is diachronic: layout, spelling changes regularly
 - Corpus is multilingual: German, French, Italian, English, Romansch etc.
 - Swiss spelling deviates from German norm
 - Many names (toponyms) that are not in OCR dictionary

Bauernhaus in der Kümme.

Photographie von *A. Blum* in Davos.

an der obern Waldgrenze (1875^m), hat der Fremdenbesuch auch im Winter zu ständiger Besiedelung geführt; im Flüelathal ist es die Paßstraße, die in dem wilden Hochgebirgsthal noch auf der „Alpenrose“ bei 1830 Metern und auf dem Tschuggen bei 1941 Metern den Menschen dem rauhen Klima trotzen läßt. Das Flüelahospiz, den größten Teil des Jahres hindurch in arktischer Umgebung (2388^m), liegt schon auf Engadinerboden. Das Dischmathal, an den Gletschern der Scaletta endigend, öde und bald des Waldes entbehrend, ist zur Winterszeit nur im vordern Thal bewohnt, und die letzten, ständig besetzten Häuser stehen schon bei 1700 Metern. Viel spärlicher bewohnt als das freundliche Sertig, kennen Flüela und Dischma kaum mehr den Ackerbau.

Bauernhaus in der Kümme

Photographie von A. Blum in Davos

an der obern Waldgrenze (1875^m), hat der Fremdenbesuch auch im Winter zu ständiger Besiedelung geführt; im Flüelathal ist es die Paßstraße, die in dem wilden Hochgebirgsthal noch auf der „Alpenrose“ bei 1830 Metern und auf dem Tschuggen bei 1941 Metern den Menschen dem rauhen Klima trotzen läßt. Das Flüelahospiz, den größten Teil des Jahres hindurch in arktischer Umgebung (2388^m), liegt schon auf Engadinerboden. Das Dischmathal, an den Gletschern der Scaletta endigend, öde und bald des Waldes entbehrend, ist zur Winterszeit nur im vordern Thal bewohnt, und die letzten, ständig besetzten Häuser stehen schon bei 1700 Metern. Viel spärlicher bewohnt als das freundliche Sertig, kennen Flüela und Dischma kaum mehr den Ackerbau.

Bauernhaus in der Kümme.

Photographie von A. Blum in Davos.

an der obern Waldgrenze (1875^m), hat der Fremdenbesuch auch im Winter zu ständiger Besiedelung geführt; im Flüelathal ist es die Paßstraße, die in dem wilden Hochgebirgsthal noch auf der „Alpenrose“ bei 1830 Metern und auf dem Tschuggen bei 1941 Metern den Menschen dem rauen Klima trotzen läßt. Das Flüelahospiz, den größten Teil des Jahres hindurch in arktischer Umgebung (2388^m), liegt schon auf Engadinerboden. Das Dischmathal, an den Gletschern der Scaletta endigend, öde und bald des Waldes entbehrend, ist zur Winterszeit nur im vordern Thal bewohnt, und die letzten, ständig besetzten Häuser stehen schon bei 1700 Metern. Viel spärlicher bewohnt als das freundliche Sertig, kennen Flüela und Dischma kaum mehr den Ackerbau.

- Compiling word lists with 19th century spelling (1500 word forms)
 - *Thal, Thür* instead of *Tal, Tür* (valley, door)
 - *passiren* instead of *passieren*
- Compiling word lists with Swiss spelling (5000 word forms)
 - ss instead of ß
- Inclusion of geographical names from Swiss Federal Office of Topography
 - 15 000 toponyms

- Performed by project members for a few books
- Time-intensive (dozens of hours per book)
- Not perfect: errors may be overlooked

- Substitution rules
- Several substitutions for a character sequence possible:
ii can be *n*, *u*, *ü*, *li* or *il*
- Most frequent¹ alternative (on a word level) is selected:
iiberein → *überein*

¹We currently use the latest Text+Berg release for our frequency counts

Our goal: Combination of two OCR systems that:

- requires little manual work
- works well with Text+Berg corpus.
- other approaches problematic:
 - dictionary-based approach suboptimal (toponyms / special terminology)
 - multilingual text

- Primary system: Abbyy Finereader 7 (AF7)
- Secondary system: OmniPage 17 (OP17)
- We train a unigram language model on the uncorrected Text+Berg corpus (created with AF7)
→ bootstrapped decision procedure
- Whenever OCR output of the two systems differs, we select the more probable (i.e. more frequent) alternative.
- In case several alternatives are equally probable, we pick AF7 output

These are not the final results!

Evaluation with manually corrected 1899 yearbook (488 pages; 220 000 tokens).

Table: Word accuracy of different OCR systems.

| system | word errors | word accuracy |
|--------------------|-------------|---------------|
| Abbyy FineReader 7 | 1260 | 99.26% |
| OmniPage 17 | 6466 | 96.21% |
| merged system | 1305 | 99.24% |

These are not the final results!

Evaluation with manually corrected 1899 yearbook (488 pages; 220 000 tokens).

Table: Word accuracy of different OCR systems.

| system | word errors | word accuracy |
|--------------------|-------------|---------------|
| Abbyy FineReader 7 | 1260 | 99.26% |
| OmniPage 17 | 6466 | 96.21% |
| merged system | 1305 | 99.24% |

Starting point for manual correction was AF7 output: evaluation is skewed in favor of AF7.

- Alignment algorithm returns 1800 differences between two systems
- In 1350 cases, AF7 alternative is chosen (trivial case)
- We manually analysed the other 450 cases:
 - 277 errors fixed
 - 82 errors added
 - 89 equally good/bad alternatives
 - **net gain: 195 improvements**, out of 1000-2000 OCR errors

Table: Examples where OmniPage is preferred over FineReader by our decision procedure.

| Abbyy FineReader | OmniPage | correct alternative in context | judgment |
|----------------------|---------------------|--|----------|
| Wunseh, | Wunsch, | entstand in unserem Herzen der Wunsch , | better |
| East | Rast | durch die Rast neu gestärkt | better |
| Übergangspunkt,. das | Übergangspunktr das | ein äußerst lohnender Übergangspunkt, das | equal |
| großen. Freude | großen, Freude | zu meiner großen Freude | equal |
| halten | hatten | Wir halten es nicht mehr aus | worse |
| là | la | c'est là le rôle principal qu'elle joue | worse |

- Combination of OCR systems improves quality
- Our decision procedure works in spite of uncorrected training data
- Selection based on unigram frequency is not optimal ($la \rightarrow la$)
Still untested: higher-order word n-grams
- Ongoing work:
 - Expansion of Text+Berg corpus
 - Further refinement of pattern-based OCR correction and OCR system combination
 - OCR error correction based on morphological analysis

Thank you for your attention!

Questions?